

Towards Symmetric Multi-threaded Optimistic Simulation Kernels

Roberto Vitali, Alessandro Pellegrini and Francesco Quaglia
Dipartimento di Ingegneria Informatica Automatica e Gestionale Antonio Ruberti
Sapienza, Università di Roma

Abstract—In this article we address the reshuffle of the design of optimistic simulation kernels in order to fit multi-core/multi-processor machines. This is done by providing a reference optimistic simulation architecture based on the symmetric multi-threaded paradigm, where each simulation kernel instance is allowed to run a dynamically changing set of worker threads that share the whole load of LPs hosted by that kernel, and that can run both application-level event handlers and kernel-level housekeeping tasks. With this organization, CPU-cores can be dynamically reassigned to the different kernels depending on fluctuations of the workload, so to maximize productivity in an orthogonal manner with respect to traditional load balancing schemes, typically employed in the context of single-threaded simulation kernels. In order to optimize efficiency and reduce wait-for-lock-release phases while synchronizing worker threads running in kernel mode, we borrow from Operating Systems’ theory by readapting the top/bottom-halves paradigm to the design of optimistic simulation systems. We also present a real implementation of our multi-threaded architecture within the ROME OpTimistic Simulator (ROOT-Sim), namely an open-source C-based simulation platform implemented according to the PDES paradigm and the optimistic synchronization approach. Experimental results for an assessment of the validity of our proposal are presented as well.

I. INTRODUCTION

In this article we focus on maximizing the productivity and the exploitation of the available computational power when running an optimistic simulation system on top of multi-core/multi-processor machines. We consider this to be a fundamental aspect nowadays, because this type of architecture has become accessible at low cost in the wide, by individuals, societies, departments and institutions. Also, the current technological trend is towards the production of chips equipped with an always increasing number of cores (many-core architectures), thus requiring optimized design/implementation approaches in order to be fruitfully exploited in the context of high performance simulations.

In the traditional approach to the design of optimistic simulation kernels, multiple LPs run within a same single-threaded simulation-kernel process (see, e.g., [1]). As a consequence, all the LPs hosted by the same kernel instance are dispatched and run on top of an individual CPU-core. Overall, parallel/distributed simulation-kernel layers commonly give control to the hosted LPs along the same thread running the CPU-scheduler (and other housekeeping tasks), according to a classical time-interleaved mode resembling what happens in traditional Operating Systems targeted at single-core machines. By this organization, the typical literature approach

aimed at achieving effective parallel/distributed simulation runs, by optimizing the exploitation of the available resources, is *load balancing*. This technique is based on migrating the application load (i.e. LPs) amongst different simulation-kernel instances (i.e. different processes). In other words, the only means to dynamically re-balance the load is to explicitly re-map the LPs across the kernels, since each kernel instance has a fixed computational power, namely one CPU-core, allocated to it.

We hereby propose a reshuffle of the internal organization of optimistic simulation kernels by presenting a symmetric multi-threaded architecture closely related to modern kernel-level Operating Systems’ technologies explicitly targeted at multi-core machines. This reshuffle supports an optimization technique orthogonal to load balancing, where the computational power (expressed in terms of CPU-cores) can be dynamically reallocated towards different active simulation-kernel instances depending on proper needs related to fluctuations in locally hosted LPs’ actual workload. Hence, we allow dynamic scale up/down in the number of worker threads belonging to each kernel instance, depending on whether locally hosted LPs (dynamically) increase/decrease their computational power demand. Overall, with our reshuffle we enable a load sharing approach, expressed in terms of dynamic redistribution of the whole simulation load across the whole set of available computational resources.

We note that the paradigm shift towards this kind of symmetric multi-threaded organization is non-trivial, since (optimistic) simulation platforms are typically expected to expose a reduced set of services (compared, e.g., to those offered by a conventional Operating System kernel), internally handled by the simulation-kernel layer via a relatively reduced set of data structures. Hence, data conflicts upon simultaneous execution in kernel mode by multiple worker threads (each running on top of a different CPU-core) may easily become a bottleneck. To address this issue, we borrow from the top/bottom-halves programming paradigm, used for handling interrupts within modern, multi-core Operating Systems, to design symmetric multi-threaded optimistic simulation kernels guaranteeing minimal length of wait-for-lock-release phases, and high scalability.

We provide a real implementation of the symmetric multi-threaded architecture within the ROOT-Sim open source optimistic simulation package [2], along with some policies for the dynamic reallocation of the computational power to

the different kernel instances. An experimental study is also presented in order to support the viability and the effectiveness of our proposal.

The remainder of this article is organized as follows. In Section II we discuss related work. The description of the symmetric multi-threaded architecture is provided in Section III. Policies for the dynamic reallocation of the computational power across the different simulation kernel instances are presented in Section IV. Section V is devoted to the experimental study.

II. RELATED WORK

The multi-threaded approach has been efficiently used in simulation platforms for separating the I/O routines from the computational ones. Similar attempts have been done in the field of HLA-based simulation platforms (see, e.g., [3]), where multi-threading has been used to implement non-blocking interoperability services across federations of simulators. The main difference from our proposal is that multi-threading has been used to implement sector-specific functionalities, while we use it as a means to overtake differentiated operations (including event processing). In addition, to the best of our knowledge, changes in the number of worker threads has never been used to perform dynamic optimizations in response to workload's variations, as instead we allow with our proposal. It has only been employed in master-slave simulation architectures to cope with dynamic increase/decrease of the amount of available resources (e.g. for simulation platforms running on top of desktop grids [4]). However, in such a context, concurrent threads operate on inherently partitioned data, while we approach multi-threading in the presence of shared (e.g. kernel-level) data structures.

When considering solutions specifically oriented to improve the performance of simulation platforms on multi-core machines, one approach related to our proposal can be found in [5]. However, this approach is targeted at a specific architecture, namely the IBM cell processor, while our proposal is general, thus being suited for differentiated multi-core platforms. Also, the work in [5] is oriented to optimize the simulation via task parallelization schemes that are orthogonal to the power reallocation scheme we present in this article.

Similar considerations can be made for other works which target simulation systems' performance improvements via the exploitation of hardware parallelism offered by GPU architectures (see, e.g., [6]). These approaches are mostly suited for data parallelism while we deal with more general schemes proper of the PDES paradigm. Also, dynamic computational power reallocation across different simulation kernel instances is not targeted by those works.

Recently, the work in [7] has presented an approach for improving the effectiveness of optimistic simulations on multi-core machines via the employment of a global schedule mechanism relying on a distributed event queue. Differently from this work, our proposal targets the traditional case of local schedule, characterized by higher scalability thanks to the avoidance of cross-kernel synchronization operations

while handling scheduling tasks. Similar considerations can be made when considering simulation architectures like ThreadedWarped [8], which uses a global priority queue. Differently from our proposal, it also uses a manager thread for event synchronizing and scheduling, which makes the architecture non-symmetric, as opposed to the symmetric approach we have devised, which does not entail any manager thread performing specialized functionalities.

Given that we provide an architectural organization which allows optimizing the use of the available computational resources in face of dynamism and fluctuations of the actual LPs' workload, our work is naturally related to all literature solutions which presented policies for load balancing in the context of either conservative (e.g. [9], [10]) or optimistic simulation (e.g. [11], [12], [13], [14]). As already hinted, the main differences between our proposal and these works are in that (A) we allow computational power reassignment, rather than workload, across the active simulation kernel instances, and (B) we rely on an innovative, multi-threading-oriented paradigm to exploit dynamically scaled up/down power available to each kernel instance. Also, our proposal can be considered as *orthogonal and complementary* to the above results, when considering that we target multi/many-core machines, while the aforementioned load balancing schemes can be used for load-redistribution on distributed memory systems (e.g. clusters).

III. THE SYMMETRIC MULTI-THREADED ARCHITECTURE

A. Handling Kernel-level Synchronization

A paradigm shift towards the design/implementation of symmetric multi-threaded optimistic simulation kernels, entailing multiple worker threads that can concurrently run any of the LPs hosted on top of the same kernel instance, needs to avoid synchronization phases while running in kernel mode to become a performance bottleneck. Specifically, while different worker threads inherently execute according to data partitioning paradigms once entered application mode (since, in accordance with the specification of the original Time Warp protocol [15], each logical process handles its own application-level data structures), care must be taken to avoid "*lock-everything effects*" when running in kernel mode. The risk for these effects is actually due to the reduced set of subsystems forming the optimistic simulation kernel (compared, e.g., to those typically included within the kernel of a general purpose Operating System), and also to the inherent strict coupling among the LPs (compared, e.g., to the typical level of coupling of different processes running on top of a conventional Operating System).

Most notably, the data structures requiring frequent updates, to be performed coherently via proper kernel-level synchronization mechanisms, are both the input and output queues of the LPs. Essentially, these data structures represent the core of cross-LP dependencies, thus involving update operations caused not only by the activities executed by the worker thread currently taking care of running the "queue-owner LP", but also by the activities carried out by worker threads taking

care of running other LPs. Synchronizing the access to these data structures via a conventional locking mechanism would give rise to scalability problems, exactly due to such a strict coupling. Further, it would give rise to critical sections whose duration would depend on the actual time-complexity of the queue-update operation.

We note that the access to the LPs' state queues (either for saving or restoring a state image) does not induce thread synchronization issues since the need for state log/recovery operations is only an indirect reflection of cross-LP coupling, caused by events scheduled across the LPs. In other words, a single worker thread is allowed to safely operate on the LP's state queue at any time, namely the worker thread that has taken care of dispatching that LP for either forward or rollback execution.

The architectural organization we propose in this paper to cope with the reduction of synchronization costs while performing housekeeping operations borrows from the design principles proper of multi-processor/multi-core Operating Systems. Specifically, any housekeeping task potentially crossing the boundaries of individual LPs' data structures is dispatched according to the same rules employed to structure modern Operating System drivers, by organizing it according to top/bottom-half activities. Hence, whenever the need for the execution of such a task arises, it (logically) takes place as an interrupt to be eventually finalized within a bottom-half module. More in details, upon the interrupt occurrence, we do not immediately finalize the task, thus not immediately locking (or waiting for the lock) on the target data structure. Instead we simply execute a light top-half module which registers the bottom-half function (and its parameters) associated with the interrupt finalization within a per-LP bottom-half queue, resembling the Linux task queue. The critical section accessing the bottom-half queue takes constant-time since each new bottom-half associated with the LP is recorder at the tail of the queue. Also, when the bottom-half tasks currently registered for a given LP are flushed, the corresponding chain of records is initially unlinked from the corresponding bottom-half queue, which is again done in constant time by unlinking the head element within the chain from its base pointer¹. Given that the access to the LP bottom-half queue represents in our architectural organization the only frequently occurring synchronization point, constant-time for the corresponding critical sections directly leads to minimizing kernel level synchronization costs.

The schematization of our proposal is presented in Figure 1. Basically, our approach can be supported by relying on a spin-lock array, named `LP_LOCKS`, having one entry for each LP hosted by the multi-threaded simulation-kernel. `LP_LOCKS[i]` is used to implement the critical section for the access to the bottom-half queue associated with the i -th LP hosted by the kernel, either for inserting a new bottom-half

¹Actual data structure updates are not performed within the critical section, but are anyhow safe since, as it will be discussed in Section III-B, for locality reasons we will allow a single worker-thread at any time to be in charge of flushing the bottom-halves of a given LP.

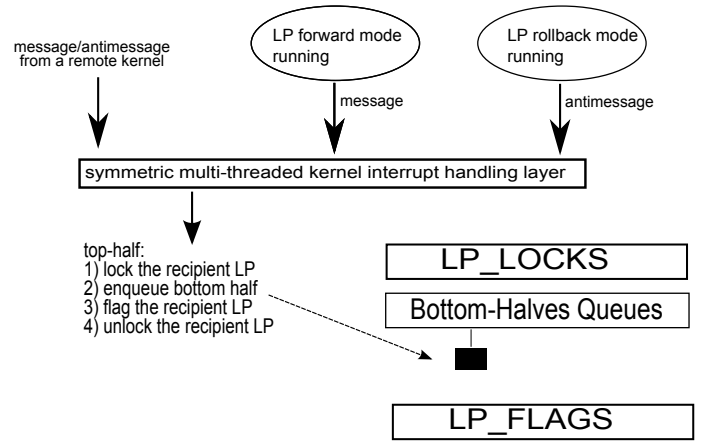


Fig. 1. Top/Bottom-Halves Architecture within the Symmetric Multi-threaded Optimistic Kernel.

task to be eventually flushed, or for taking care of unlinking the current chain, in order to flush the pending bottom-halves.

Let us now depict when (logical) interrupts to be handled via this type of organization occur. Basically, an interrupt occurs as soon as any worker thread currently active within the symmetric multi-threaded kernel becomes aware of a new message/antimessage destined to the i -th locally hosted LP. In such a case, the worker thread needs to access the i -th bottom-half queue within a critical section that performs the insertion of the corresponding message/antimessage delivery task, as explained above. To provide additional details, awareness by a worker thread of a new message/antimessage destined to a locally hosted LP arises in three different circumstances:

- (i) The worker thread is currently running the locally hosted LP_j in forward mode, and this LP produces a new event to be scheduled for the locally hosted LP_i . Thus the worker thread enters kernel mode for actuating the delivery of the corresponding message to LP_i 's input-queue. (Note that j might be equal to i , hence giving rise to the case where sender and receiver coincide.)
- (ii) The worker thread is currently running the locally hosted LP_j in rollback mode (hence it is performing kernel level housekeeping operations associated with revealed causality errors), which gives rise to the production of an antimessage destined to LP_i , which requires access to LP_i 's input queue for annihilating the original message. (Again we might have $j = i$.)
- (iii) The message passing layer notifies the worker thread (e.g. via an explicit message receive operation executed by this thread according to a traditional polling scheme) about a new message/antimessage incoming from some remote kernel instance.

As shown in Figure 1, we logically mark all the above three circumstances as interrupts, which will be treated homogeneously, and whose associated message/antimessage delivery operation will be finalized via the bottom-half mechanism.

We note that spin-locks may anyhow exhibit non-minimal costs since they require the corresponding operations to be

performed via sequences of atomic instructions (e.g. via the LOCK prefix for the IA-32 instruction set). Additionally, since they are shared and accessed by different threads, cross-cache invalidation effects can be induced as soon as one worker thread gains control on the spin-lock. To reduce these effects, we have devised the presence of an additional array of flags LP_FLAGS (see again Figure 1), where LP_FLAGS[i] indicates whether the corresponding bottom-half queue, namely the one associated with the i -th locally hosted LP, is not empty. Actually, LP_FLAGS[i] gets updated within a critical section protected by LP_LOCKS[i], either when a new bottom-half is inserted within the corresponding queue (in this case the flag is raised), or when the queue is flushed (in this case the flag is reset). However, LP_FLAGS[i] is also accessed before trying to lock the bottom-half queue in order to avoid spin-lock operations in all the cases where the queue would reveal empty once accessed within the critical section leading to flush operations. The exact scheme looks therefore as follows:

```

TOP-HALF:                                BOTTOM-HALF:
lock(&LP_LOCKS[i]);                       if (LP_FLAGS[i]){
<log bottom-half>;                          if (try_lock(&LP_LOCKS[i])){
LP_FLAGS[i] = TRUE;                          <unlink bottom-halves>;
unlock(&LP_LOCKS[i]);                          LP_FLAGS[i] = FALSE;
                                                unlock(&LP_LOCKS[i]);
                                                <perform bottom-halves>;
                                                }
}

```

Being LP_FLAGS[i] checked non-atomically wrt lock acquisition when attempting to perform bottom-halves, we might experience false negatives in case the top-half finalizes the insertion of the bottom-half task concurrently with the check. However, this does not represent a safety problem since the flag will be rechecked periodically in subsequent attempts to flush the corresponding bottom-half queue, thus eventually falling within the case where the bottom-half queue is correctly reflected into the state of the input queue of the destination LP. Such a reflection might therefore experience only a delay, which resembles delays introduced by traditional single-threaded kernels while reflecting the content of cross-kernel messages into the system state, which is typically affected by the polling period according to which the messaging layer is accessed for acquiring not yet delivered messages. Further, as hinted in footnote 1, a single worker thread at a time will be allowed to manage flush operations for a given LP, hence no false positives will ever be experienced.

As a last note, messages/antimessages whose deliveries are still pending, being them recorded as tasks to be finalized within bottom-half queues, represent a kind of in-transit data, whose timestamp needs to be accounted for when computing the GVT value.

B. Tackling Locality Issues

Given that all the worker threads associated with the same simulation kernel instance operate within the same address space, the symmetric multi-threaded kernel allows virtual addresses related to both application and kernel level data structures (associated with whichever LP) to be, in principle,

accessible by any worker thread. However, such a level of sharing would cause frequent invalidation/refill of, e.g., the top-level private caches of individual cores, even when entailing processor affinity schemes involving the worker threads. As an example, data structures associated with an LP that has been lastly accessed by a given worker thread would be flushed by the corresponding private caching system upon the first write access by a different worker thread.

Overall, while developing a symmetric multi-threaded optimistic simulation-kernel a core additional issue to address is related to maintaining an adequate level of locality, so to avoid harming caching performance. In order to cope with this issue, we devise the adoption of affinity mechanisms such that a worker thread belonging to a given simulation-kernel instance is not allowed to run every LP hosted by that kernel. Instead, it takes care of running a subset of these LPs, which are currently selected as being affine to the worker thread. In other words, we devise the use of temporary binding mechanisms associating a subset of the locally hosted LPs to a specific worker thread, which is therefore the only thread taking care of running these LPs during a specific wall-clock-time window. We note that this approach resembles what is done by the scheduler of Linux kernel 2.6, where a temporary binding of active processes/threads to a specific CPU-core is supported for both (a) locality and (b) reduction of the CPU scheduling cost.

Overall, within the affinity scheme, each worker thread is in charge of:

- (i) Flushing the bottom-half queues associated with its affine LPs, which is executed periodically according to a traditional polling approach.
- (ii) Dispatching its affine LPs for execution in time interleaved mode.

We note anyway that the binding of a specific LP to a worker thread is not meant to be fixed, but can change overtime, also in relation to variations of the amount of worker threads activated within a given symmetric optimistic kernel instance. The policy according to which the locally hosted LPs are reassigned to the worker threads will be discussed in Section IV, together with the performance model we use to reallocate the computational power (and hence CPU-cores) to the different symmetric simulation-kernel instances.

As an additional note in relation to locality, we also devise proper memory layout mechanisms in order to reduce the false cache sharing problem for kernel-level data structures. As an example, the entries of both the LP_LOCKS and LP_FLAGS arrays, which represent frequently accessed synchronization data structures, can be memory bind to different cache lines, which can be easily achieved by exploiting, e.g., the posix_memalign API plus padding schemes. The same approach can be taken for the meta-data associated with each single LP hosted by each instance of the symmetric multi-threaded kernel, so that once an LP is bind to a given worker thread, cache interference due to accesses to meta-data does not arise.

IV. COMPUTATIONAL POWER REALLOCATION POLICIES

The symmetric multi-threaded kernel allows scaling up/down the amount of per-kernel worker threads without any change in the internal operating mode. This allows for dynamically reallocating the computational power (in terms of CPU-cores) to the active kernel instances depending on fluctuations of the workload and efficiency variations within the optimistic simulation run. In this section we first provide an approach for reallocating the CPU-cores to the active kernels. Then we address the issue of (temporarily) binding the LPs hosted by a given simulation kernel instance to specific worker threads.

A. Dynamical Assignment of CPU-Cores to Kernels

Let us denote with C_{tot} the amount of available CPU-cores, and with $K_{tot} < C_{tot}$ the number of active symmetric multi-threaded kernel instances (the case $K_{tot} = C_{tot}$ trivially boils down to the traditional scenario where each kernel instance is allowed to run on a single CPU-core, hence in the typical single-threaded mode). Our first objective is to determine the amount of CPU-cores C_i (with $1 \leq i \leq K_{tot}$) to be assigned to kernel instance K_i for a given wall-clock-time window, so to improve resource exploitation for fruitful processing activities.

In our proposal, the re-evaluation of C_i values can be carried out periodically, for example upon computing a new GVT value or after a set of subsequent GVT computations. This also allows to exploit a set of metrics characterizing the parallel simulation run, as an example in terms of determination of the event rate (committed events per wall-clock-time unit) achieved by each of the symmetric multi-threaded kernel instances. We denote the event rate achieved by kernel K_i as evr_i . This quantity is a measure for the fruitful (non-rolled back) amount of simulation work carried out by each kernel instance. In an ideal scenario where the efficiency is maximized (i.e. where the undone computation is negligible), each symmetric multi-threaded kernel instance K_i should use an amount of computational power that suffices to execute exactly evr_i events per wall-clock-time unit. In fact, an excess of computational power could lead to over-optimism and hence to rolled back computations, thus moving the run-time dynamics far from the above depicted ideal case. So the idea behind the determination of C_i values is to dynamically assign an amount of CPU-cores to kernel K_i which is proportional to the actual computation requirements of K_i for the achievement of its relative event rate, compared to the one by the other kernels. Actually, to also take care of the real CPU requirements on a given kernel instance (so to also take into account possible variance of the event granularity across the LPs hosted by different kernel instances), which is the indicator of the real usage of computational power for committing the events, the evr_i metric can be refined by weighting it via the average CPU time required for processing the events on a specific kernel K_i , which we denote as Δ_i . Hence we express the weighted event rate as $wavr_i = evr_i \times \Delta_i$.

In other words, $wavr_i$ values observed during the last wall-clock-time period express the relative CPU requirements of

each kernel instance in order to carry out productive simulation work, in relation to the activities of the other kernels and the outcoming synchronization dynamics. Hence, assigning a computational power proportional to the relative weighted event rate would tend to lead to the situation where each kernel instance advances its LPs in simulation time in a “synchronization suited” manner according to what the other kernels are able to do on their own. This part of the dynamic reallocation scheme would therefore tend to avoid significant presence of overoptimistic kernel instances during the various phases of the run.

It is anyway typical that performance can be further enhanced even in cases where the efficiency is already maximized (or optimized), for example by further reassigning the computational power depending on the real weight of the workload associated with the hosted LPs. As an example, for loosely synchronized models we may have two or more groups of LPs that do not interact, or stop interacting during the run (hence eventually not directly impacting synchronization and efficiency), exhibiting different speed of advancement in simulation time due to, e.g., different weights of the corresponding events in terms of CPU requirements. In such a case, the completion of the simulation would be delayed by the slowest group. Therefore, within the dynamic scheme for resource assignment, an increase of computational power should also be envisaged for all those kernel instances exhibiting larger CPU requirements to advance in simulation time. To this end we include in our scheme the parameter $wcta_i$, which indicates the wall-clock-time required by kernel K_i to advance a single simulation time unit. The usage of this parameter within the dynamic reallocation scheme would tend to complement the above described one by further attempting to align the advancement of the different symmetric multi-threaded kernel instances in simulation time while the run proceeds.

Finally, the amount of cores C_i to be assigned to kernel K_i should anyway be bounded by the maximum degree of parallelism that can be accomplished by K_i , which is a function of the amount of locally hosted LPs. In fact, each LP is an intrinsically sequential entity, which is not further parallelized, thus not being allowed to simultaneously use multiple CPU-cores for its execution.

Overall, we devise the following rules for dynamically defining the amount of CPU-cores to be reassigned to each kernel K_i in order to optimize the usage of the available computational power:

- 1) For each kernel K_i the parameter $\alpha_i = \frac{wavr_i}{\sum_{j=1}^{K_{tot}} wavr_j}$ is computed.
- 2) A first calculation of C_i is then performed as $C_i = \lfloor \alpha_i \times C_{tot} \rfloor$.
- 3) For each kernel instance K_i for which the condition $C_i \geq numLP_i$ is verified (where $numLP_k$ identifies the number of LPs hosted by K_i), then C_i is definitively set to $numLP_i$. In fact, additional CPU-cores could not be effectively exploited for parallelization of the locally hosted LPs.

- 4) At this point, there could be some CPU-cores left to be assigned, which we decide to assign on the basis of (A) the request for allocation remainder of kernel K_i , namely $r_i = [(\alpha_i \times C_{tot}) - C_i]$ and (B) the parameter $wcta_i$. In particular, we order the kernels for which the finalization of C_i values still needs to be performed (so the ones already finalized in point 3 are excluded) according to decreasing values of the product $r_i \times wcta_i$, and we assign the remaining CPU-cores according to a round-robin rule following the priority defined by such an ordering.

Each of the above steps is an implementation of the rationales discussed above in terms of suited CPU-core assignment vs specific performance aspects.

B. Binding LPs to Worker Threads

As pointed out earlier, a given set of LPs hosted by K_i gets temporarily bind to a specific worker thread acting within the kernel, which is in charge of performing bottom-half operations related to the LPs in the set, and to schedule them for event processing according to some priority scheme (e.g. Lowest-Timestamp-First). Once the new value for C_i gets defined upon reallocating the computational power, a policy is required to determine which LPs are bind to a specific worker thread. To achieve a binding that allows balancing the whole workload related to local LPs onto the whole set of worker threads, we have devised the below policy. For the j -th LP hosted by kernel K_i , which we refer to as LP_i^j , we compute the total amount of CPU-time required for committing its events during the last observation period (e.g. the last GVT cycle). We refer to this metric as cpu_i^j .

The maximum cpu_i^j value across all the locally hosted LPs represents in our scheme a reference knapsack, and the corresponding LP_i^j is assigned to a given worker thread. Then we exploit the greedy approximation approach proposed by George Dantzig in [16] which allows a maximum “overflow” of about 30% over the reference knapsack, in order to build the other knapsacks of LPs (hence knapsacks characterized by sums of cpu_i^x values) to be assigned to the remaining worker threads. We do this by actually applying a variant of the original scheme, where the knapsacks are filled according to a round-robin approach. The procedure is then iterated until no more LP needs to be further bind to any worker thread.

V. EXPERIMENTAL STUDY

A. Test-bed Platform

We have implemented the proposed symmetric multi-threaded optimistic kernel architecture within ROOT-Sim, which is an open source C/MPI-based simulation package targeted at POSIX systems [2], which implements a general-purpose parallel/distributed simulation environment relying on the optimistic synchronization paradigm.

ROOT-Sim offers a very simple programming model based on the classical notion of simulation-event handlers (both for processing events and for accessing a committed and globally

consistent state image upon GVT calculations), to be implemented according to the ANSI-C standard, and transparently supports all the services required to parallelize the execution. It also offers a set of optimized protocols aimed at minimizing the run-time overhead by the platform, thus allowing for high performance and scalability.

Among the main features offered by ROOT-Sim we can mention completely transparent recoverability of the state of the LPs achieved through proper hooking of dynamic memory allocation/release [17], plus ad-hoc code instrumentation schemes that allow incremental determination of dirty state portions [18] and that, ultimately, allow dynamical switch between different state log/restore schemes depending on the proper dynamics of the application layer [19].

The single threaded version of ROOT-Sim also offers innovative transparent supports for LP migration and load balancing [20], which will be considered as a reference for the assessment of the currently presented symmetric multi-threaded version in terms of ability to exploit the computational resources offered by a multi-core machine when the actual simulation workload dynamically varies over time.

Future steps ahead in the development of ROOT-Sim definitely entail the integration of the currently presented symmetric multi-threaded architecture with the aforementioned LP migration subsystem, currently supported only when running in single-threaded mode. This will ultimately provide an environment where the orthogonal capabilities offered by the symmetric multi-threaded paradigm (in terms of dynamic reassignment of the computational power to different kernel instances) and the traditional migration approach (in terms of ability to move individual LPs across different kernel instances) get ultimately combined.

Integration of the multi-threaded approach within ROOT-Sim has been based on `pthread` technology, and on the reorganization of the kernel level data structures in order to (i) provide per-thread private data, and (ii) cache aligned kernel-level memory buffers so to avoid false cache sharing across the worker threads within the same symmetric multi-threaded kernel instance. The latter target has been achieved by exploiting the `posix_memalign` API, plus the usage of proper padding schemes allowing cache alignment for sequences of records, such as arrays of values. As for the accesses to the MPI layer, used to transfer messages across different kernel instances, in our architecture they can be symmetrically issued by any of the worker thread operating within a given kernel instance. Given that the MPI layer does not natively support multi-threading, we have included a wrapper that synchronizes these accesses transparently towards the worker threads via the embedding of critical sections protected by spin-locks.

As far as GVT computation and fossil collection are concerned, we have implemented a symmetric scheme where upon a new GVT computation, all the worker threads operating within a same kernel instance run a race. The race winner actually computes the local reduction and interacts with the master kernel in order to determine the globally reduced value representing the new GVT. However, once defined the

new GVT value, all the worker threads operating within the same kernel instance are allowed to perform fossil collection operations in parallel. Each of these threads takes care of fossil collecting the obsolete information associated with its affine LPs.

Finally, the hardware architecture used for testing our proposal is a 64-bit NUMA machine, namely an HP Proliant server, equipped with four 2GHz AMD Opteron 6128 processors and 64GB of RAM. Each processor has 8 CPU-cores (for a total of 32 CPU-cores) that share a 10MB L3 cache (5118KB per each 4-cores set), and each core has a 512KB private L2 cache. The operating system is 64-bit Debian 6, with Linux kernel version 2.6.32.5. The compiling and linking tools used are gcc 4.4.5 and binutils (as and ld) 2.20.0.

B. Application Benchmarks

In order to evaluate different aspects of the proposed symmetric multi-threaded architecture, we have conducted experiments on two different application benchmarks, namely *PCS (Personal Communication System)* and *Traffic*, which are hereby described. The first one has been configured in order to provide a constant workload across all the LPs during the whole simulation run. This has been done in order to measure the actual overhead of the symmetric multi-threaded architecture, while not taking advantages from its ability to reallocate CPU-cores just given the constancy of the workload. The second application benchmark provides instead a highly dynamic workload that varies over time across the involved LPs. This type of benchmark has been used in order to assess the goodness of the symmetric multi-threaded architecture in terms of its ability to reallocate the computational power depending on the actual needs.

1) *The PCS Benchmark*: this application benchmark implements a simulation model of wireless communication systems adhering to GSM technology, where communication channels are modeled in a high fidelity fashion via explicit simulation of power regulation/usage and interference/fading phenomena on the basis of the current state of the corresponding cell. The power regulation model has been implemented according to the results in [21].

Upon the start of a call destined to a mobile device currently hosted by a given wireless cell, a call-setup record is instantiated via dynamically-allocated data structures, which gets linked to a list of already active records within that same cell. Each record gets released when the corresponding call ends or is handed-off towards an adjacent cell. In the latter case, a similar call-setup procedure is executed at the destination cell. Upon call-setup, power regulation is performed, which involves scanning the aforementioned list of records for computing the minimum transmission power allowing the current call-setup to achieve the threshold-level SIR value. Data structures keeping track of fading coefficients are also updated while scanning the list, according to a meteorological model defining climatic conditions (and related variations). The climatic model accounts for variations of the climatic

conditions (e.g. the current wind speed) with a minimum time granularity of ten seconds.

This simulation model has been developed for execution on top of ROOT-Sim in a way that each LP models a single wireless cell. Hence, the event-handler callback involves the update of individual cells' states, and cross-LP events are essentially related to hand-offs between different cells.

To evaluate the overhead due to the symmetric multi-threaded architecture, when compared to the classical case of single-threaded optimistic kernel, we have performed a set of experiments where each wireless cell sustains the same workload of incoming calls, hence we are in a balanced scenario not requiring dynamical reallocation of the computational power, which is instead a main target of the symmetric multi-threaded organization. The call inter-arrival time is exponentially distributed, and the average call duration is set to 2 min. The expected rate for call inter-arrival has been set to achieve channel utilization factor on the order of 30%, while the residence time of an active device within a cell has a mean value of 5 min and follows the exponential distribution. For the above scenario, we have run experiments with 1024 wireless cells, modeled as hexagons covering a square region, each one managing 1000 wireless channels.

We have measured the cumulated event rate (expressed as the amount of cumulated committed events vs wall-clock-time) for different configurations of the symmetric multi-threaded kernel, comparing it with the one achievable when running the same ROOT-Sim package in single-threaded mode. In particular, executions with 4, 8, 16 and 32 symmetric multi-threaded kernels (each one starting with 8, 4, 2 and 1 worker thread, respectively) have been carried out. Also, in order to assess the effects of the symmetric multi-threaded organization, we additionally report statistics related to typical run-time parameters characterizing optimistic simulation runs, such as the rollback frequency and the rollback length.

2) *The Traffic Benchmark*: this benchmark application simulates a complex highway system (at a single car granularity), where the topology is a generic graph, where nodes represent cities or junctions and edges represent the actual highways. Every node is described in terms of car inter-arrival time and car leaving probability, while edges are described in terms of their length.

At startup phase, the simulation model is asked to distribute the highway's topology on a given number of LPs. Every LP therefore handles the simulation of a node or a portion of a segment, the length of which depends on the total highway's length and the number of available LPs.

Cars enter the system according to an Erlang probability distribution, with a mean interarrival time specified (for each node) in the topology configuration file. They can join the highway starting from cities/junctions only, and are later directed towards highway segments with a uniform probability. Whenever a car is received, it is enqueued in the LP's list of traversing cars, and its speed (for the particular LP it is entering in) is determined according to a Gaussian probability distribution, the mean and the variance of which are specified

at startup time. Then, the model computes the time the car will need to traverse the node, adding traffic slowdowns which are again computed according to a Gaussian distribution. In particular, the probability of finding a traffic jam is a function of the number of cars which are currently passing through the node.

Accidents are derived according to a probability function as well. In particular, they are more likely to occur when the amount of cars traversing an LP is about half of the cars which can be hosted altogether. In fact, if few cars are in, accidents are less frequent. Similarly, if there are many, the traffic factor produces a speed slowdown, entailing the probability of an accident to occur to be reduced. Therefore, the model discretizes a Normal distribution, computing the Cumulative Density Function in a contour defined as *cars in the node* $\pm \frac{1}{2}$, having as the mean half of the total number of cars which are at the current moment in the system, and as variance a factor which can be specified at startup. The total number of cars which can be hosted by an LP is computed according to the actual length of the simulated road, which is determined when the model is initialized. When an accident occurs, the cars are not allowed to leave the LP, until the road is freed. The duration of an accident phase is determined according to a Gaussian distribution, the mean and the variance of which are again specified at startup.

In our execution, we have simulated the whole Italian highway network on top of 1024 LPs. We have discarded the highways segments in the islands in order to simulate an undirected connected graph, which allows to have the actual workload migrating overall the highway. The topology has been derived from [22], and the traffic parameters have been tuned according to the measurements provided in [23]. The average speed has been set to 110 Km/h, with a variance of 20 Km/h, and accident durations have been set to 1 hour, with 30 minutes variance. This model has provided results which are statistically close to the real measurements provided in [24].

We consider this second application benchmark to be significant for showing how our proposed symmetric multi-threaded architecture is able to capture unbalance in the load, and react via computational power reallocation across the active kernels in order to drive the system back into an evenly-distributed workload processing scenario, which would lead to enhanced fruitful exploitation of the computational resources.

For this benchmark application we still report the event rate, this time comparing it with both the one achieved when considering the classical single-threaded execution mode of ROOT-Sim, and the one achievable when activating within such a single-threaded mode the load balancing mechanisms described in [20]. We recall again that these load balancing facilities are in principle orthogonal to the facilities offered by the symmetric multi-threaded organization, since their target is the move of LPs across the kernels, not the reassignment of the computational power to the multi-threaded kernels. Anyway, we feel that taking load balancing facilities properly offered by the single-threaded version of ROOT-Sim into account in

the comparison provides a relevant reference for assessing the potential offered by the symmetric multi-threaded organization in terms of its ability to fruitfully exploit the available computational power with dynamic workloads.

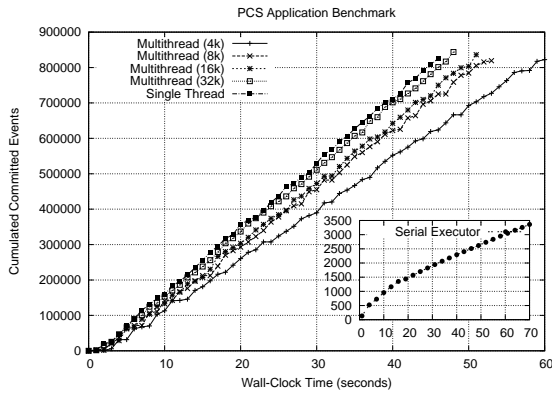
C. Results

In Figure 2 we show the experimental results that have been obtained for the PCS application benchmark. We recall that this benchmark exhibits balanced workload during the whole run, hence it is suited for assessing the overhead of the symmetric multi-threaded architecture when considering that its capabilities to redistribute the computational power across the different kernel instances are not actually exploited. By the curves related to the cumulated committed events (where all the samples have been obtained as the average over 10 runs all done with different pseudo-random seeds) we see that, unless for the case of 4 multi-threaded kernels (each running 8 worker threads), the additional latency for reaching the completion of the simulation run, compared to the traditional case of single-threaded kernel, is no more than 13%. This is an indication of limited performance intrusiveness by the top/bottom-half architecture while managing the input/output queues of the LPs, as well as limited performance intrusiveness by other synchronization mechanisms (e.g. for the access to the MPI layer), at least when the scale of the multi-threaded configuration is bounded by the value 4. On the other hand, the multi-threaded configuration with 4 kernels and 8 worker threads for each kernel exhibits a non-minimal overhead (on the order of 25%).

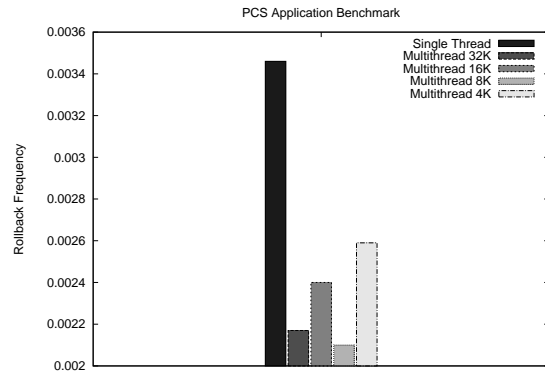
We note that, although the workload is constant, small fluctuations due to the probability distribution ruling the generation of the events can arise, which are therefore captured by our symmetric multi-threaded architecture. Nevertheless, this sensibility enhances the reassignment overhead, as long as the time spent in this operation is not rewarded by the new worker threads' configuration. This is more likely to occur exactly when the average number of worker threads per kernel instance gets increased. In addition, we note that these data have been achieved by considering a model with medium-to-fine event granularity, on the order of 30/40 microsecs, thus further supporting the viability of our proposal, since applications exhibiting coarser-grained events would absorb better the actual overhead of the multi-threaded architecture.

Also, we note that the parallel runs provide a super-scalar speedup with respect to the serial executor (based on the calendar-queue scheduler), which indicates that the experimentation has been carried out when considering competitive parallel runs.

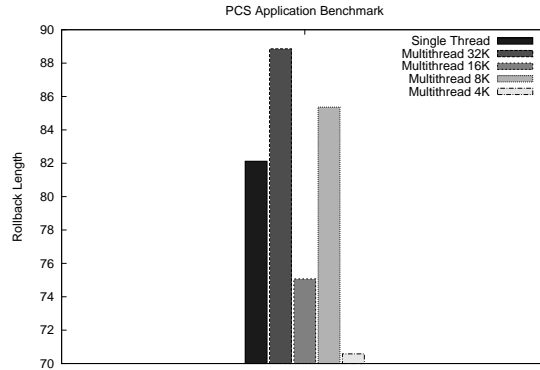
For completeness, in Figure 2 we also report the observed values for rollback frequency and rollback length for the PCS application benchmark. By these data we can observe how the symmetric multi-threaded kernel tends to exhibit a slightly throttled execution profile, compared to the single-thread case. In particular, we note a clear reduction of the rollback frequency, with a less significant increase of the rollback length. For the case of the symmetric multi-threaded



(a) Cumulated Committed Events

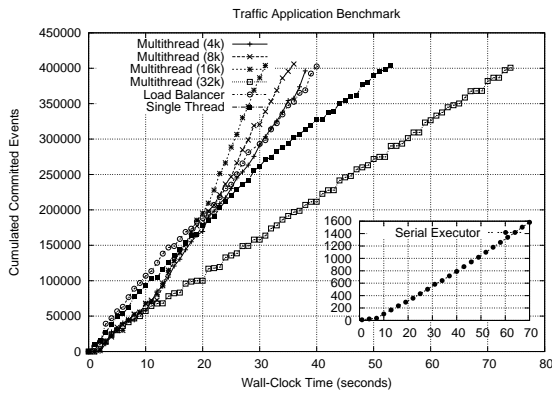


(b) Rollback Frequency

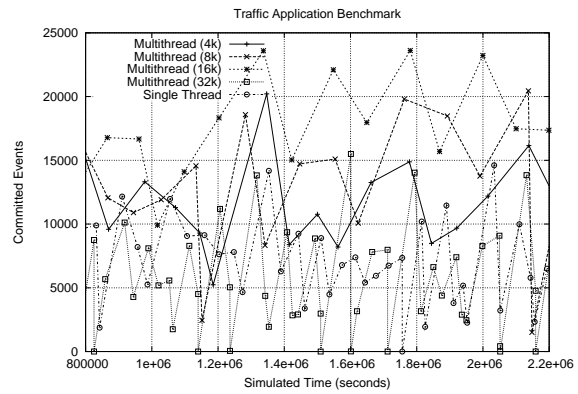


(c) Rollback Length

Fig. 2. Results for the PCS Application Benchmark.



(a) Cumulated Committed Events



(b) Punctual Event Rate

Fig. 3. Results for the Traffic Application Benchmark.

configuration with 4 kernels, such a throttling is an expression of the above noted overhead, which leads to less favorable run-time dynamics.

In Figure 3 we report the results for the case of the Traffic application benchmark. This time we have compared the cumulated event rate by our symmetric multi-threaded architecture with a classical single-threaded organization, a serial execution of the same application-level software running

on top of a calendar-queue scheduler, and also results of the load balancing architecture based on the migration approach presented in [20]. Again, the parallel approaches provide a super-scalar speedup. The multi-threaded versions of the simulation kernel provide a speedup wrt the single-threaded one, which ranges in between 35% (for the 4 kernels configuration) and 73% (for the 16 kernels configuration).

As for the 32 multi-threaded kernels execution, we note that

the speed down is in the order of 37%. This is related to the fact that in this configuration no actual power reallocation is possible on the 32-core server machine that has been employed (in fact, each simulation kernel must have at least one worker thread in order to proceed in the simulation). Therefore, we are again simply measuring the symmetric multi-threaded architecture pure overhead.

The last comparison shown by the plots is the one wrt the traditional load balancer. Although we note that the load balancer configuration provides a speedup in the order of 30% wrt the single-threaded approach, it's throughput is comparable with the 4 kernels multi-threaded configuration, while the 8 and 16 kernels configurations of the multi-threaded architecture are still 30% faster than the traditional load balancer configuration.

As a final note, always in Figure 3 we report the punctual variation of the event rate over time for the case of a single run (hence not mediated over different runs). These data show how the dynamical reassignment of resources, depending on fluctuations of the workload, leads the symmetric multi-threaded architecture to provide punctual improvements in the amount of committed events per wall-clock-time unit, which are quantified by these plots and are then ultimately reflected in the above discussed performance improvements.

VI. CONCLUSIONS AND FUTURE WORK

In this article we have presented the design and implementation of a symmetric multi-threaded optimistic simulation kernel targeted at multi-core/multi-processor machines, where, similarly to what happens in multi-core oriented Operating Systems in terms of process management, multiple threads operate symmetrically in order to sustain the whole workload associated with the LPs hosted by a kernel instance. This type of organization allows to transparently scale up/down the amount of worker threads operating within a same instance of the optimistic simulation kernel. Hence, it allows for dynamic reassignment of the computational power, namely CPU-cores, to the different kernel instances involved within the optimistic run, depending on variations of the workload associated with the hosted LPs. Policies suited for the reassignment have been also presented, and the whole system has been tested with different application benchmarks.

As future work we plan to integrate the symmetric multi-threaded architecture with traditional load balancing facilities, thus eventually allowing the optimistic kernel to reconfigure its run-time behavior towards an optimal use of the available resources by jointly exploiting CPU-core reassignment facilities and traditional LP migration schemes.

REFERENCES

- [1] D. E. Martin, T. J. McBrayer, and P. A. Wilsey, "WARPED: A Time Warp simulation kernel for analysis and application development," in *Proceedings of the 29th Hawaii International Conference on System Sciences (HICSS), Volume 1: Software Technology and Architecture*. IEEE Computer Society, 1996, p. 383.
- [2] F. Quaglia, A. Pellegrini, and R. Vitali, "ROOT-Sim: The ROME Optimistic Simulator: <http://www.dis.uniroma1.it/~hpdc/root-sim/>," Oct. 2011.
- [3] L. Mellon and D. West, "Architectural optimizations to advanced distributed simulation," in *Proceedings of Winter Simulation Conference*, 1995, pp. 634–641.
- [4] A. Park and R. Fujimoto, "Optimistic parallel simulation over public resource-computing infrastructures and desktop grids," in *Proceedings of the 12th IEEE/ACM International Symposium on Distributed Simulation and Real Time Applications (DS-RT)*, 2008, pp. 149–156.
- [5] Q. Liu and G. Wainer, "Multicore acceleration of discrete event system specification systems," *SIMULATION*, 2011. [Online]. Available: <http://sim.sagepub.com/content/early/2011/06/28/0037549711412237.abstract>
- [6] T. Hamada and K. Nitadori, "190 tflops astrophysical n-body simulation on a cluster of gpus," in *Proceedings of the Conference on High Performance Computing Networking, Storage and Analysis*. IEEE Computer Society, 2010, pp. 1–9.
- [7] L. li Chen, Y. shuai Lu, Y. ping Yao, S. liang Peng, and L. da Wu, "A well-balanced Time Warp system on multi-core environments," in *Proceedings of the 25th ACM/IEEE Workshop on Principles of Advanced and Distributed Simulation (PADS)*, 2011, pp. 154–162.
- [8] R. J. Miller, *Optimistic Parallel Discrete Event Simulation on a Beowulf Cluster of Multi-core Machines*. Cincinatti University: Master Dissertation, 2010.
- [9] A. Boukerche and S. K. Das, "Dynamic load balancing strategies for conservative parallel simulations," in *Proceedings of the 11th ACM/IEEE International Workshop on Parallel and Distributed Simulation (PADS)*, 1997, pp. 20–28.
- [10] G. D'Angelo and M. Bracuto, "Distributed simulation of large-scale and detailed models," *International Journal of Simulation and Process Modelling (IJSPM)*, vol. 5, no. 2, pp. 120–131, 2009.
- [11] D. W. Glazer and C. Tropper, "On process migration and load balancing in Time Warp," *IEEE Transactions on Parallel and Distributed Systems*, vol. 4, no. 3, pp. 318–327, 1993.
- [12] C. D. Carothers and R. Fujimoto, "Efficient execution of Time Warp programs on heterogeneous, NOW platforms," *IEEE Transactions on Parallel and Distributed Systems*, vol. 11, no. 3, pp. 299–317, 2000.
- [13] P. L. Reiher and D. Jefferson, "Virtual time based dynamic load management in the Time Warp operating system," *Transactions of the Society for Computer Simulation*, vol. 7, pp. 103–111, 1990.
- [14] S. Meraji, W. Zhang, and C. Tropper, "A multi-state q-learning approach for the dynamic load balancing of Time Warp," in *Proceedings of the 24th ACM/IEEE International Workshop on Principles of Advanced and Distributed Simulation (PADS)*, 2010, pp. 1–8.
- [15] D. R. Jefferson, "Virtual Time," *ACM Transactions on Programming Languages and System*, vol. 7, no. 3, pp. 404–425, Jul. 1985.
- [16] G. B. Dantzig, "Discrete-variable extremum problems," *Operational Research*, no. 5, pp. –, 1957.
- [17] R. Toccaceli and F. Quaglia, "DyMeLoR: Dynamic memory logger and restorer library for optimistic simulation objects with generic memory layout," in *Proceedings of the 22nd ACM/IEEE International Workshop on Principles of Advanced and Distributed Simulation (PADS)*. IEEE Computer Society, 2008, pp. 163–172.
- [18] A. Pellegrini, R. Vitali, and F. Quaglia, "Di-DyMeLoR: Logging only dirty chunks for efficient management of dynamic memory based optimistic simulation objects," in *Proceedings of the 23rd ACM/IEEE International Workshop on Principles of Advanced and Distributed Simulation (PADS)*. IEEE Computer Society, 2009, pp. 45–53.
- [19] R. Vitali, A. Pellegrini, and F. Quaglia, "Autonomic log/restore for advanced optimistic simulation systems," in *Proceedings of the 18th Annual IEEE/ACM International Symposium on Modeling, Analysis and Simulation of Computer and Telecommunication Systems (MASCOTS)*. Los Alamitos, CA, USA: IEEE Computer Society, 2010, pp. 319–327.
- [20] S. Peluso, D. Didona, and F. Quaglia, "Application transparent migration of simulation objects with generic memory layout," in *Proceedings of the 25th ACM/IEEE International Workshop on Principles of Advanced and Distributed Simulation (PADS)*. IEEE Computer Society, 2011, pp. 169–177.
- [21] S. Kandukuri and S. Boyd, "Optimal power control in interference-limited fading wireless channels with outage-probability specifications," *IEEE Transactions on Wireless Communications*, vol. 1, no. 1, pp. 46–55, 2002.
- [22] "Atlante stradale italia," <http://www.automap.it/>.
- [23] http://www.autostrade.it/studi/studi_traffico.html.
- [24] "Aci - dati e statistiche," <http://www.aci.it/?id=54>.